

Universidade Federal Fluminense
Instituto de Ciências Exatas - ICEx



Computação de Alto Desempenho 2

Paralelização da solução para equação de Schrödinger

Bianca Maia

23 de julho de 2022

Resumo

Ao reconhecer a necessidade de aumentar o poder computacional para o interesse acadêmico em pesquisas científicas, descobertas e também otimizações na operação de uma grande quantidade de dados, um novo tipo de máquinas foi desenvolvido, estas são chamadas de supercomputadores. A computação de alto desempenho (HPC) é o uso ativo desses supercomputadores e clusters para a operação de tarefas que exigem grande poder de processamento e problemas em que a solução numérica é muito complicada de se obter. A maioria desses problemas numéricos contém seus dados dentro de matrizes, portanto, uma forma de melhorá-los com computação paralela é por unidades de processamento gráfico (GPU) que possuem uma arquitetura paralela massiva que o torna extremamente eficaz nas operações com matrizes.

Um exemplo bem conhecido de um problema complicado, em relação ao tempo de processamento, são as equações diferenciais parciais. Este trabalho tem como foco resolver uma equação de Schrödinger pelo método das diferenças finitas, otimizando seu tempo de execução com CUDA. Apresenta também técnicas de melhoria de código e análise de flags.

Palavras-chaves: HPC, computação paralela, programação paralela com GPU, CUDA

Abstract

By recognizing the urge to increase computational power for academic interest in scientific research, discoveries and also optimizations in operation some large amount of data, a new type of machines had been developed, these are called supercomputers. High-performance computing (HPC) is the active use of these supercomputers and cluster in order to operating tasks that require a large processing power and problems in which the numerical solution are too complicated to obtain. Most of these numerical problems contains its data inside matrices, therefore a way of improving it with parallel computing it's by graphics processing units (GPU) which have a massive parallel architecture that makes it extremely effective in matrices operations.

A well-known example of a complicated problem, in respect of processing time, is the partial differential equations. This works focus on solving a Schrödinger equation by the finite difference method, optimizing its execution time with CUDA. It also presents techniques of code improvement and flags analysis.

Keywords: HPC, parallel computing, GPU parallel programming, CUDA

Conteúdo

1	Introdução	5
2	Objetivo	5
3	Metodologia	6
3.1	Discretização para a equação de Schrödinger	6
4	Implementação (fluxograma)	7
5	Benchmark (serial)	8
5.1	Lab 107	8
5.1.1	Análise para o compilador gnu	9
5.1.2	Análise para o compilador ifort	10
5.2	SequanaX	11
5.2.1	Análise para o compilador gnu	12
5.2.2	Análise para o compilador ifort	13
5.3	Comentários	14
6	Profile (análise dos gargalhos)	15
6.1	Profile sem flags	15
6.2	Profile com a melhor flag	18
6.3	Análise de resultados	19
7	Técnicas de otimização de software	20
8	Comparação do programa base com o otimizado	21
8.1	Profile com otimizações nível código (sem flags)	21
8.2	Análise de resultados	24
9	Implementação do programa com CUDA	26
9.1	Estratégia de paralelização	26
9.2	Profile	26
10	Benchmark (paralelo)	27
10.1	Tesla K40t (LNCC)	28
10.1.1	Análise de flags	30
10.2	Tesla V100 (LNCC)	36
10.2.1	Análise de flags	38
10.3	Benchmark com as melhores flags (Tesla V100)	43
11	Comparação com a implementação em OpenMP	45
11.1	Considerações sobre o método das diferenças finitas	46
12	Validação de Resultados	47
13	Conclusões	49

14 Apêndice	50
14.1 Informações sobre as CPUs	50
14.1.1 Laboratório 107C (UFF)	50
14.1.2 B107 (LNCC)	50
14.1.3 SequanaX (LNCC)	51
14.2 Informações sobre as GPUs	52
14.2.1 Tesla K40t	52
14.2.2 Tesla V100	53

1 Introdução

Considere a equação de Schrödinger dependente do tempo em duas dimensões que descreve como o estado quântico de um sistema físico evolui com o tempo

$$-\frac{\hbar}{2m}\nabla^2\psi(\vec{r}) + V(\vec{r})\psi(\vec{r}) = E\psi(\vec{r}) \quad (1)$$

onde $\psi(\vec{r}) = \psi(x, y, t)$, $V(\vec{r}) = V(x, y)$ e $\nabla^2 = \frac{\partial^2\psi}{\partial x^2} + \frac{\partial^2\psi}{\partial y^2}$.

O objetivo do algoritmo é resolver um problema onde um elétron está confinado dentro de um potencial quadrático com a sua posição inicial dada pela gaussiana:

$$\psi(x, y, t = 0) = e^{ik_{0x}x}e^{ik_{0y}y}\exp\left(-\frac{(x-x_0)^2}{2\sigma_0^2}\right)\exp\left(-\frac{(y-y_0)^2}{2\sigma_0^2}\right) \quad (2)$$

O potencial utilizado para a resolução computacional desse problema foi:

$$V(x, y) = 0.01x^2 + 0.07y^2$$

O programa para a solução dessa equação com essas condições foi inicialmente desenvolvido na linguagem C em serial, dispensando quaisquer preocupações em relação ao seu tempo de execução. Posteriormente, esse foi otimizado por programação paralela utilizando o openMP com o objetivo de reduzir o tempo de compilação.

2 Objetivo

Resolver uma equação diferencial parcial computacionalmente pode ser custoso em termos de tempo de compilação dependendo das condições iniciais determinadas para o problema.

O principal objetivo do projeto é reduzir o tempo de execução do programa com o uso da programação paralela, sem perder as características do problema da equação de Schrödinger.

Na finalidade de atingir essa otimização, serão analisados os desempenhos do programa serial para um conjunto de flags referentes à dois compiladores em dois hardwares diferentes. As melhores flags serão aproveitadas para uma análise do programa em paralelo para um conjunto de threads (referentes à cada máquina).

Na necessidade de otimizações à nível código no serial, um profiling apontará quais dependências do programa serial mais consomem desempenho para aplicar as boas práticas de programação eficientemente.

Finalmente, com as melhores flags e as melhores threads, o programa paralelizado com openMP poderá obter uma otimização de tempo mais agressiva.

3 Metodologia

A metodologia usada para resolver o problema de equação diferencial parcial foi o **método das diferenças finitas** [1].

3.1 Discretização para a equação de Schrödinger

Reescrevendo a equação (1) como

$$\psi(x, y, t) = U(t)\psi(x, y, t = 0) \quad (3)$$

onde $U(t) = e^{-i\tilde{H}t}$ e $\tilde{H} = -\frac{\partial^2\psi}{\partial x^2} + \frac{\partial^2\psi}{\partial y^2} + V(x, y)$. Temos que $U(t)$ é um operador que desenvolve a função de onda por um determinado tempo t e \tilde{H} é o operador Hamiltoniano.

Assumindo que o pacote de onda avançará de acordo com um intervalo temporal definido por Δt

$$\psi_{i,j}^{n+1} = U(\Delta t)\psi_{i,j}^n \quad (4)$$

o sobrescrito denota o tempo $t = n\Delta t$ e os sobrescritos denotam variáveis espaciais $x = i\Delta x$ e $y = j\Delta y$.

Dada a inversa que retorna à solução no intervalo anterior

$$\psi^{n-1} = U^{-1}(\Delta t)\psi^n \quad (5)$$

e a relação de diferença entre ψ^{n+1} e ψ^{n-1}

$$\psi^{n+1} = \psi^{n-1} + [e^{-i\tilde{H}\Delta t} - e^{+i\tilde{H}\Delta t}]\psi^n \quad (6)$$

da expansão de Taylor,

$$\frac{\partial^2\psi}{\partial x^2} \approx \frac{-1}{2}[\psi_{i+1,j}^n + \psi_{i-1,j}^n - 2\psi_{i,j}^n] \quad (7)$$

logo,

$$\psi_{i,j}^{n+1} = \psi_{i,j}^{n-1} - 2i[(4\alpha + \frac{1}{2}\Delta t V_{i,j})\psi_{i,j}^n - \alpha(\psi_{i+1,j}^n + \psi_{i-1,j}^n + \psi_{i,j+1}^n + \psi_{i,j-1}^n)] \quad (8)$$

onde $\alpha = \Delta t/2(\Delta x)^2$.

Para providenciar a conservação numérica da função densidade de probabilidade, o método da diferenças finitas resolve a parte real e a parte imaginária [2]. Portanto,

$$R_{i,j}^n = R_{i,j}^{n-1} + 2[(4\alpha + \frac{1}{2}\Delta t V_{i,j})I_{i,j}^n - \alpha(I_{i+1,j}^n + I_{i-1,j}^n + I_{i,j+1}^n + I_{i,j-1}^n)]$$

$$I_{i,j}^n = I_{i,j}^{n-1} - 2[(4\alpha + \frac{1}{2}\Delta t V_{i,j})R_{i,j}^n + \alpha(R_{i+1,j}^n + R_{i-1,j}^n + R_{i,j+1}^n + R_{i,j-1}^n)]$$

4 Implementação (fluxograma)

O fluxograma do código utilizado para resolver o problema da equação de Schrödinger está disposto abaixo.

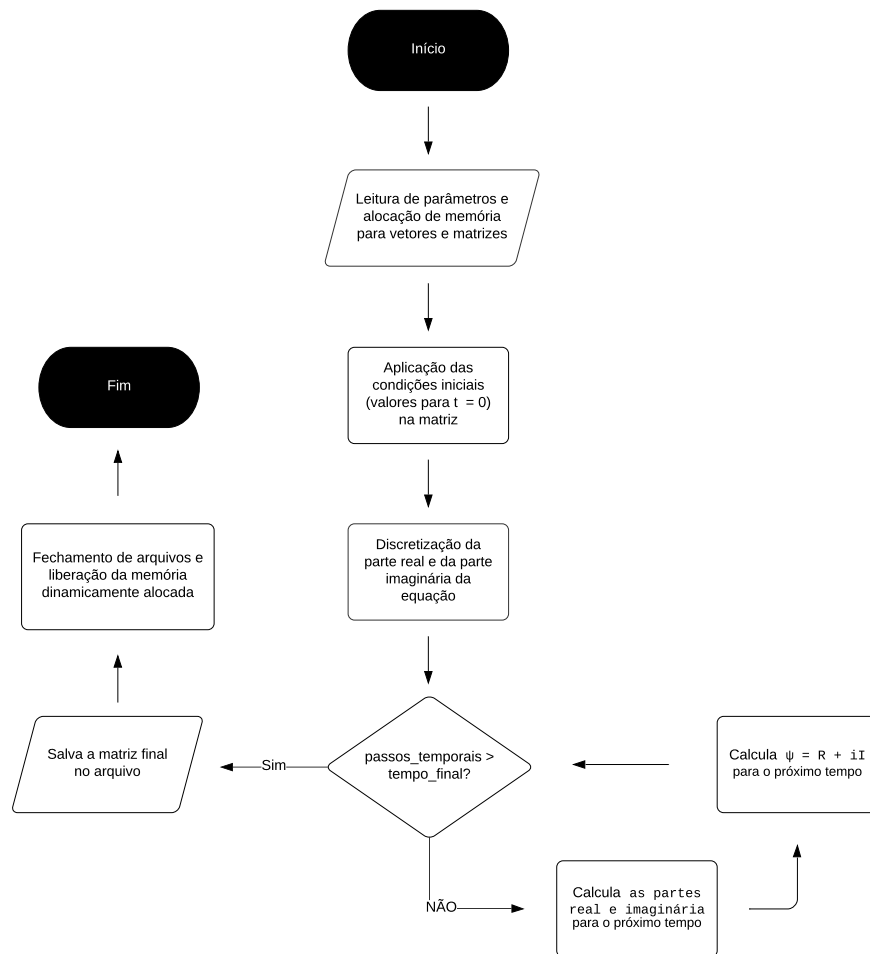


Figura 1: Fluxograma do programa para resolver a equação de Schrödinger com o método das diferenças finitas.

5 Benchmark (serial)

O benchmark foi realizado apenas para o Lab 107 e o SequanaX, devido à indisponibilidade do B710 durante o período em que essas análises foram realizadas.

5.1 Lab 107

Nesse computador as flags utilizadas para o compilador gnu foram:

1. -OX
2. -OX -fexpensive-optimizations -m64 -foptimize-register-move -funroll-loops -ffast-math
3. -OX -mtune=native -march=native -fexpensive-optimizations -m64 -foptimize-register-move -funroll-loops -ffast-math
4. -OX -mtune=corei7-avx -march=corei7-avx -fexpensive-optimizations -m64 -foptimize-register-move -funroll-loops -ffast-math

e para o compilador ifort:

1. -OX
2. -OX -w -mp1 -zero -xHOST
3. -OX -w -mp1 -zero -xHOST -fast=2

onde as otimizações -OX compreendem a variação entre -O0, -O1, -O2 e -O3.

5.1.1 Análise para o compilador gnu

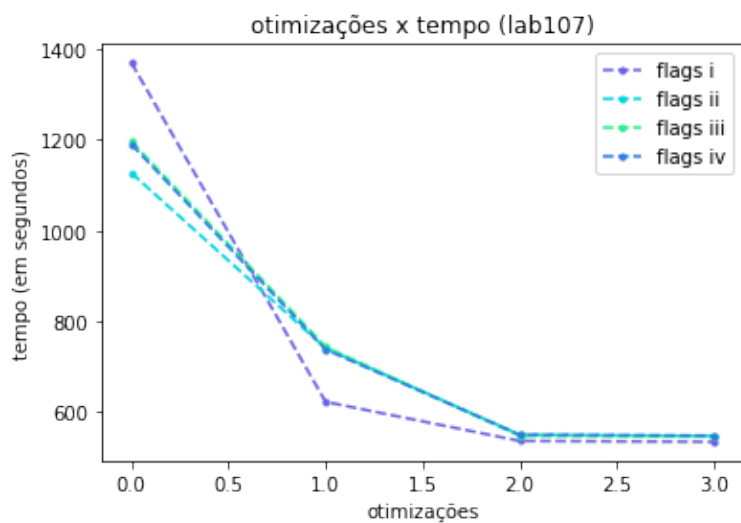


Figura 2: Gráfico para os tempos em função das otimizações com as flags do compilador gnu no computador do Lab 107.

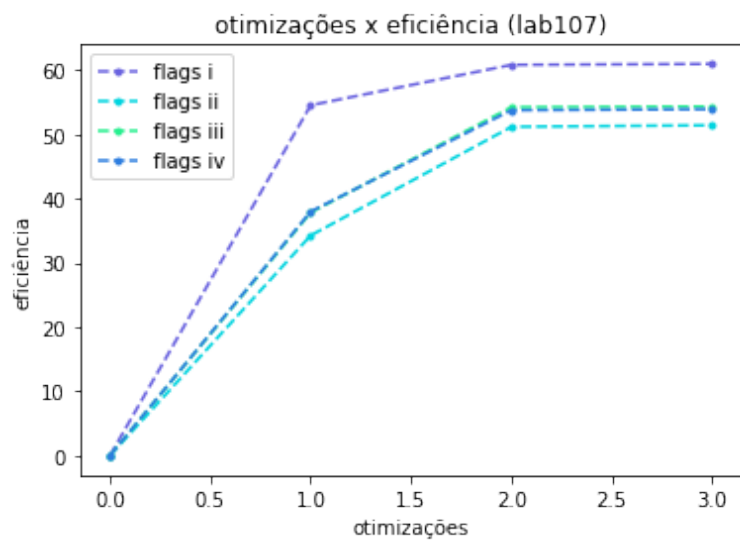


Figura 3: A eficiência utilizando as flags do compilador gnu no Lab 107.

É facilmente observável que a melhor flag gnu para esse computador

é a -O2 do conjunto de flags i. Embora a -O3 do mesmo conjunto seja uma concorrente em eficiência, pelo gráfico do tempo vemos que a outra supera essa última.

5.1.2 Análise para o compilador ifort

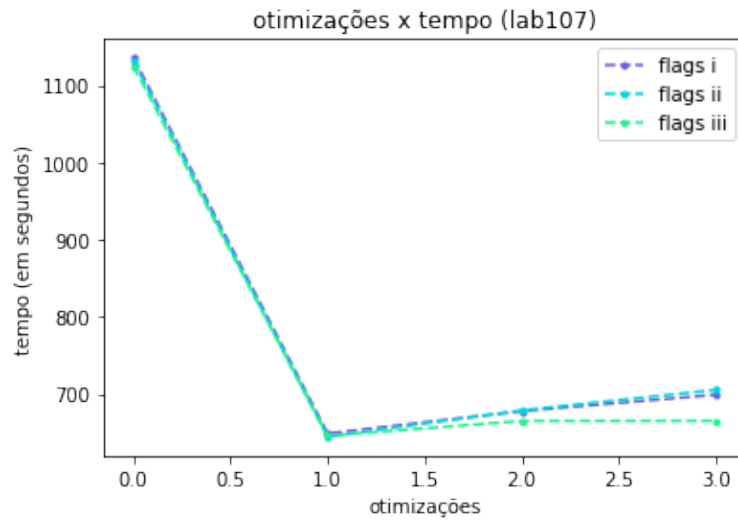


Figura 4: O tempo para o compilador ifort no Lab 107.

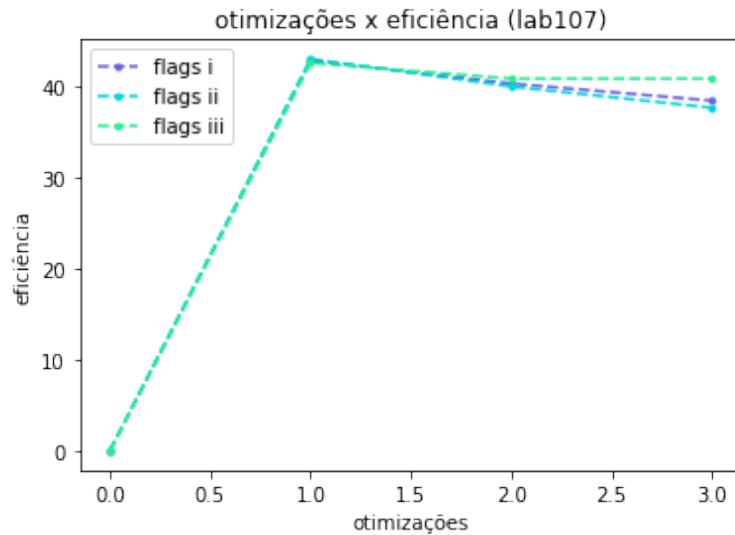


Figura 5: Gráfico da eficiência para o compilador ifort no Lab 107.

Nesse compilador, é imediato que a melhor flag é a -O1 para quaisquer um dos conjuntos de flags utilizados.

5.2 SequanaX

Nesse computador as flags utilizadas para o compilador gnu foram:

1. -OX
2. -OX -fexpensive-optimizations -m64 -foptimize-register-move -funroll-loops -ffast-math
3. -OX -mtune=native -march=native -fexpensive-optimizations -m64 -foptimize-register-move -funroll-loops -ffast-math
4. -OX -mtune=core-avx2 -march=core-avx2 -fexpensive-optimizations -m64 -foptimize-register-move -funroll-loops -ffast-math

e para o compilador ifort:

1. -OX
2. -OX -w -mp1 -zero -xHOST
3. -OX -w -mp1 -zero -xHOST -fast=2

onde as otimizações -OX compreendem a variação entre -O0, -O1, -O2 e -O3.

5.2.1 Análise para o compilador gnu

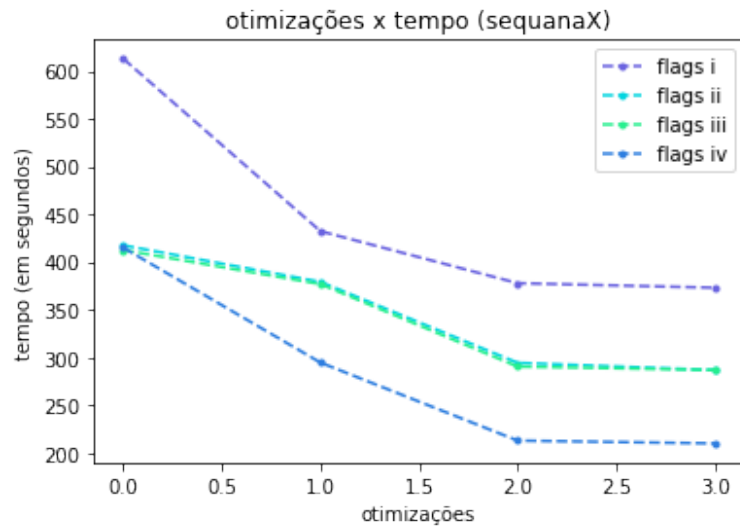


Figura 6: O tempo com as flags gnu para o SequanaX.

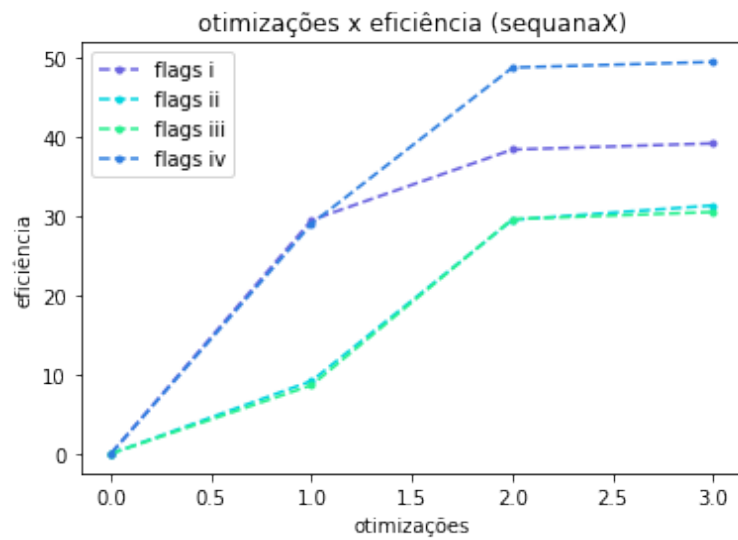


Figura 7: Gráfico da eficiência para o compilador gnu no SequanaX.

Diferentemente da situação anterior para o compilador gnu, as melhores flags são aquelas acompanhadas do último conjunto de flags que é acompanhando

de flags referentes à arquitetura do computador. Nesse conjunto, tanto a -O2 quanto a -O3 são as mais eficientes que todas as outras, embora a -O3 ligeiramente apresente um desempenho superior à -O2.

5.2.2 Análise para o compilador ifort

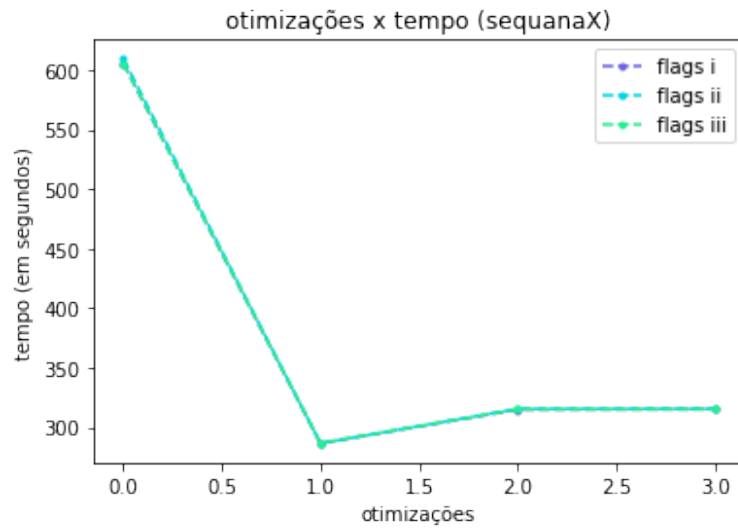


Figura 8: Tempo para o compilador ifort no SequanaX.

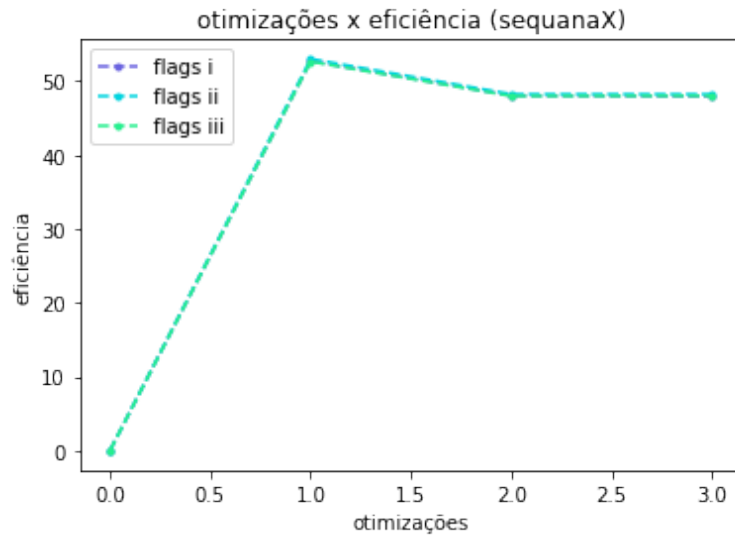


Figura 9: A eficiência para as flags ifort no SequanaX.

Nesse último caso, as curvas são praticamente concorrentes. As melhores flags são aquelas com otimização -O1 para todos os conjuntos utilizados nesse computadores.

5.3 Comentários

Para o compilador gnu, a melhor flag na máquina do Lab 107 é a:

```
-lm -O2
```

Por outro lado, no SequanaX o impacto das otimizações de arquitetura são mais intensos e portanto a melhor flag é:

```
-lm -O3 -mtune=core-avx2 -march=core-avx2
-fexpensive-optimizations -m64 -foptimize-register-move
-funroll-loops -ffast-math
```

Para o compilador ifort, em ambos os computadores as melhores flags são aquelas acompanhadas da otimização -O1 independente do conjunto de flags. Não há critérios rigorosos para decidir qual conjunto utilizar na análise do código em paralelo, dado que a diferença de tempo entre esses é ignorável.

6 Profile (análise dos gargalhos)

6.1 Profile sem flags

O profile foi executado para o programa serial do projeto com a matriz reduzida. O tempo de execução do programa está logo abaixo, antes dos resultados obtidos pelo profiling.

```
1 real 35m15,495s
2 user 35m5,720s
3 sys 0m4,657s
```

```
1 Flat profile:
2
3 Each sample counts as 0.01 seconds.
4 % cumulative self self total
5 time seconds seconds calls ms/call ms/call name
6 100.30 2117.86 2117.86 600000 3.53 3.53
7 finite_difference
8 0.00 2117.90 0.04
9 0.00 2117.92 0.02 1 20.07 20.07 main
10 0.00 2117.92 0.00 3 0.00 0.00 wave_packet
11 0.00 2117.92 0.00 1 0.00 0.00 allocate
12 allocate_complex
13 0.00 2117.92 0.00 1 0.00 0.00 potential
14
15 % the percentage of the total running time of the
16 time program used by this function.
17
18 cumulative a running sum of the number of seconds accounted
19 seconds for by this function and those listed above it.
20
21 self the number of seconds accounted for by this
22 seconds function alone. This is the major sort for this
23 listing.
24
25 calls the number of times this function was invoked, if
26 this function is profiled, else blank.
27
28 self the average number of milliseconds spent in this
29 ms/call function per call, if this function is profiled,
30 else blank.
31
32 total the average number of milliseconds spent in this
33 ms/call function and its descendents per call, if this
34 function is profiled, else blank.
35
36 name the name of the function. This is the minor sort
37 for this listing. The index shows the location of
38 the function in the gprof listing. If the index is
39 in parenthesis it shows where it would appear in
40 the gprof listing if it were to be printed.
41
42 Copyright (C) 2012-2020 Free Software Foundation, Inc.
```

```

43 Copying and distribution of this file, with or without modification
44 ,
45 are permitted in any medium without royalty provided the copyright
46 notice and this notice are preserved.
47
48 Call graph (explanation follows)
49
50
51 granularity: each sample hit covers 2 byte(s) for 0.00% of 2117.92
52 seconds
53
54 index % time self children called name
55 [1] 100.0 0.04 2117.88 <spontaneous>
56 2117.86 0.00 600000/600000 main [1]
57 [2] 0.02 0.00 1/1 wave_packet [3]
58 0.00 0.00 3/3 allocate [4]
59 0.00 0.00 1/1 allocate_complex
60 [5] 0.00 0.00 1/1 potential [6]
61 -----
62 2117.86 0.00 600000/600000 main [1]
63 [2] 100.0 2117.86 0.00 600000 finite_difference [2]
64 -----
65 0.02 0.00 1/1 main [1]
66 [3] 0.0 0.02 0.00 1 wave_packet [3]
67 -----
68 0.00 0.00 3/3 main [1]
69 [4] 0.0 0.00 0.00 3 allocate [4]
70 -----
71 0.00 0.00 1/1 main [1]
72 [5] 0.0 0.00 0.00 1 allocate_complex [5]
73 -----
74 0.00 0.00 1/1 main [1]
75 [6] 0.0 0.00 0.00 1 potential [6]
76 -----
77
78 This table describes the call tree of the program, and was sorted
79 by
80 the total amount of time spent in each function and its children.
81
82 Each entry in this table consists of several lines. The line with
83 the
84 index number at the left hand margin lists the current function.
85 The lines above it list the functions that called this function,
86 and the lines below it list the functions this one called.
87 This line lists:
88 index A unique number given to each element of the table.
89 Index numbers are sorted numerically.
90 The index number is printed next to every function name so
91 it is easier to look up where the function is in the table.
92
93 % time This is the percentage of the `total' time that was
94 spent
95 in this function and its children. Note that due to

```



```

93 different viewpoints, functions excluded by options, etc,
94 these numbers will NOT add up to 100%.
95
96     self This is the total amount of time spent in this function.
97
98     children This is the total amount of time propagated into this
99     function by its children.
100
101     called This is the number of times the function was called.
102     If the function called itself recursively, the number
103     only includes non-recursive calls, and is followed by
104     a '+' and the number of recursive calls.
105
106     name The name of the current function. The index number is
107     printed after it. If the function is a member of a
108     cycle, the cycle number is printed between the
109     function's name and the index number.
110
111
112 For the function's parents, the fields have the following meanings
113 :
114
115     self This is the amount of time that was propagated directly
116     from the function into this parent.
117
118     children This is the amount of time that was propagated from
119     the function's children into this parent.
120
121     called This is the number of times this parent called the
122     function '/' the total number of times the function
123     was called. Recursive calls to the function are not
124     included in the number after the '/'.
125
126     name This is the name of the parent. The parent's index
127     number is printed after it. If the parent is a
128     member of a cycle, the cycle number is printed between
129     the name and the index number.
130
131 If the parents of the function cannot be determined, the word
132 '<spontaneous>' is printed in the 'name' field, and all the other
133 fields are blank.
134
135 For the function's children, the fields have the following
136 meanings:
137
138     self This is the amount of time that was propagated directly
139     from the child into the function.
140
141     children This is the amount of time that was propagated from
142     the
143     child's children to the function.
144
145     called This is the number of times the function called
146     this child '/' the total number of times the child
147     was called. Recursive calls by the child are not
148     listed in the number after the '/'.

```

```

147     name This is the name of the child. The child's index
148     number is printed after it. If the child is a
149     member of a cycle, the cycle number is printed
150     between the name and the index number.
151
152 If there are any cycles (circles) in the call graph, there is an
153 entry for the cycle-as-a-whole. This entry shows who called the
154 cycle (as parents) and the members of the cycle (as children.)
155 The '+' recursive calls entry shows the number of function calls
    that
156 were internal to the cycle, and the calls entry for each member
    shows,
157 for that member, how many times it was called from other members
    of
158 the cycle.
159
160
161 Copyright (C) 2012-2020 Free Software Foundation, Inc.
162
163 Copying and distribution of this file, with or without modification
    ,
164 are permitted in any medium without royalty provided the copyright
165 notice and this notice are preserved.
166
167
168 Index by function name
169
170     [4] allocate                [2] finite_difference    [6]
        potential
171     [5] allocate_complex        [1] main                [3]
        wave_packet

```

6.2 Profile com a melhor flag

A matriz é a mesma do profile anterior.

```

1 real 10m32,015s
2 user 10m30,145s
3 sys 0m1,260s

1 Flat profile:
2
3 Each sample counts as 0.01 seconds.
4   % cumulative self      self      total
5   time  seconds seconds    calls  Ts/call  Ts/call  name
6 100.40    633.97   633.97
    finite_difference
7   0.00    633.98    0.01
                                           wave_packet
8
9   %           the percentage of the total running time of the
10  time        program used by this function.
11
12 cumulative a running sum of the number of seconds accounted
13 seconds    for by this function and those listed above it.
14
15 self       the number of seconds accounted for by this

```

```

16 seconds      function alone.  This is the major sort for this
17               listing.
18
19 calls         the number of times this function was invoked, if
20               this function is profiled, else blank.
21
22 self          the average number of milliseconds spent in this
23 ms/call       function per call, if this function is profiled,
24               else blank.
25
26 total        the average number of milliseconds spent in this
27 ms/call       function and its descendents per call, if this
28               function is profiled, else blank.
29
30 name          the name of the function.  This is the minor sort
31               for this listing. The index shows the location of
32               the function in the gprof listing. If the index is
33               in parenthesis it shows where it would appear in
34               the gprof listing if it were to be printed.
35
36
37 Copyright (C) 2012-2020 Free Software Foundation, Inc.
38
39 Copying and distribution of this file, with or without modification
40 ,
41 are permitted in any medium without royalty provided the copyright
   notice and this notice are preserved.

```

6.3 Análise de resultados

A função que mais consome processamento é a função responsável por resolver a equação utilizando o método das diferenças finitas.

O profile sem flags não resolve o problema de tempo de execução, é mais eficiente com as flags que provocam uma melhoria de aproximadamente 20min.

7 Técnicas de otimização de software

Como demonstrado no profile, a função responsável pelo consumo de tempo é a função do método das diferenças finitas. Nada é possível fazer a respeito da otimização do código dessa função, no entanto ainda é viável conferir outras funções essenciais para a resolução do problema, a função do potencial e a do pacote de onda.

```
1 void potential(double **potential){
2     int i, j;
3     double n = 110;
4
5     for(i = 0; i < size; i++){
6         for(j = 0; j < size; j++){
7             potential[i][j] = 0.07*pow(i, 2) + 0.01*pow(j, 2);
8         }
9     }
10 }
11
12 void wave_packet(double complex **p, double **im, double **r,
13     double k0x, double k0y, double x0, double y0, double sigma){
14     int x, y, it = 0;
15
16     for(x = 0; x < size; x++){
17         for(y = 0; y < size; y++){
18             r[x][y] = cos(k0x*x + k0y*y)*exp(-(pow(x - x0, 2) + pow
19 (y - y0, 2))/(2*pow(sigma, 2)));
20             im[x][y] = sin(k0x*
21 x + k0y*y)*exp(-(pow(x - x0, 2) + pow(y - y0, 2))/(2*pow(sigma,
22 2)));
23             p[x][y] = r[x][y] + I*im[x][y];
24         }
25     }
26 }
```

As otimizações aplicáveis nessas funções são simples, consistem em reduzir o uso de funções pertencentes à biblioteca math.h quando deseja-se calcular o quadrado de determinada variável.

```
1 void potential(double **potential){
2     int i, j;
3     double n = 110;
4
5     for(i = 0; i < size; i++){
6         for(j = 0; j < size; j++){
7             potential[i][j] = 0.07*i*i + 0.01*j*j;
8         }
9     }
10 }
11
12 void wave_packet(double complex **p, double **im, double **r,
13     double k0x, double k0y, double x0, double y0, double sigma){
14     int x, y, it = 0;
15
16     for(x = 0; x < size; x++){
17         for(y = 0; y < size; y++){
18             r[x][y] = cos(k0x*x + k0y*y)*exp(-(((x - x0)*(x - x0))
19 + ((y - y0)*(y - y0)))/(2*(sigma*sigma)));
20         }
21     }
22 }
```

```

18         im[x][y] = sin(k0x*x + k0y*y)*exp(-(((x - x0)*(x - x0))
19         + ((y - y0)*(y - y0)))/(2*(sigma*sigma)));
20         p[x][y] = r[x][y] + I*im[x][y];
21     }
22 }
23 }

```

8 Comparação do programa base com o otimizado

Considera-se novamente a matriz reduzida das análises anteriores.

8.1 Profile com otimizações nível código (sem flags)

```

1 real 29m15,086s
2 user 29m12,914s
3 sys 0m0,744s

```

```

1 Flat profile:
2
3 Each sample counts as 0.01 seconds.
4
5 % cumulative self self total
6 time seconds seconds calls ms/call ms/call name
7 96.86 1699.38 1699.38 600000 2.83 2.83
8   finite_difference
9   0.00 1699.40 0.02
10  0.00 1699.40 0.00 3 0.00 0.00 main
11  0.00 1699.40 0.00 1 0.00 0.00 allocate
12   allocate_complex
13  0.00 1699.40 0.00 1 0.00 0.00 potential
14  0.00 1699.40 0.00 1 0.00 0.00 wave_packet
15
16 %
17 time the percentage of the total running time of the
18      program used by this function.
19
20 cumulative a running sum of the number of seconds accounted
21 seconds for by this function and those listed above it.
22
23 self the number of seconds accounted for by this
24 seconds function alone. This is the major sort for this
25 listing.
26
27 calls the number of times this function was invoked, if
28      this function is profiled, else blank.
29
30 self the average number of milliseconds spent in this
31 ms/call function per call, if this function is profiled,
32      else blank.
33
34 total the average number of milliseconds spent in this
35 ms/call function and its descendents per call, if this
36      function is profiled, else blank.

```

```

33
34 name          the name of the function.  This is the minor sort
35                for this listing.  The index shows the location of
36                the function in the gprof listing.  If the index is
37                in parenthesis it shows where it would appear in
38                the gprof listing if it were to be printed.
39
40
41 Copyright (C) 2012-2020 Free Software Foundation, Inc.
42
43 Copying and distribution of this file, with or without modification
44 ,
45 are permitted in any medium without royalty provided the copyright
46 notice and this notice are preserved.
47
48                Call graph (explanation follows)
49
50
51 granularity: each sample hit covers 2 byte(s) for 0.00% of 1699.40
52                seconds
53
54 index % time    self  children    called    name
55 [1]    100.0    0.02 1699.38
56                1699.38    0.00 600000/600000    main [1]
57                [2]
58                0.00    0.00    3/3    allocate [3]
59                0.00    0.00    1/1    allocate_complex
60                [4]
61                0.00    0.00    1/1    potential [5]
62                0.00    0.00    1/1    wave_packet [6]
63 -----
64                1699.38    0.00 600000/600000    main [1]
65 [2]    100.0 1699.38    0.00 600000    finite_difference [2]
66 -----
67                0.00    0.00    3/3    main [1]
68 [3]    0.0    0.00    0.00    3    allocate [3]
69 -----
70                0.00    0.00    1/1    main [1]
71 [4]    0.0    0.00    0.00    1    allocate_complex [4]
72 -----
73                0.00    0.00    1/1    main [1]
74 [5]    0.0    0.00    0.00    1    potential [5]
75 -----
76                0.00    0.00    1/1    main [1]
77 [6]    0.0    0.00    0.00    1    wave_packet [6]
78 -----
79
80 This table describes the call tree of the program, and was sorted
81 by
82 the total amount of time spent in each function and its children.
83
84 Each entry in this table consists of several lines.  The line with
85 the
86 index number at the left hand margin lists the current function.
87 The lines above it list the functions that called this function,

```

```

84 and the lines below it list the functions this one called.
85 This line lists:
86     index A unique number given to each element of the table.
87     Index numbers are sorted numerically.
88     The index number is printed next to every function name so
89     it is easier to look up where the function is in the table.
90
91     % time This is the percentage of the `total' time that was
    spent
92     in this function and its children. Note that due to
93     different viewpoints, functions excluded by options, etc,
94     these numbers will NOT add up to 100%.
95
96     self This is the total amount of time spent in this function.
97
98     children This is the total amount of time propagated into this
99     function by its children.
100
101     called This is the number of times the function was called.
102     If the function called itself recursively, the number
103     only includes non-recursive calls, and is followed by
104     a `+' and the number of recursive calls.
105
106     name The name of the current function. The index number is
107     printed after it. If the function is a member of a
108     cycle, the cycle number is printed between the
109     function's name and the index number.
110
111
112 For the function's parents, the fields have the following meanings
    :
113
114     self This is the amount of time that was propagated directly
115     from the function into this parent.
116
117     children This is the amount of time that was propagated from
118     the function's children into this parent.
119
120     called This is the number of times this parent called the
121     function `/' the total number of times the function
122     was called. Recursive calls to the function are not
123     included in the number after the `/'.
124
125     name This is the name of the parent. The parent's index
126     number is printed after it. If the parent is a
127     member of a cycle, the cycle number is printed between
128     the name and the index number.
129
130 If the parents of the function cannot be determined, the word
131 `' is printed in the `name' field, and all the other
132 fields are blank.
133
134 For the function's children, the fields have the following
    meanings:
135
136     self This is the amount of time that was propagated directly
137     from the child into the function.

```

```

138
139     children This is the amount of time that was propagated from
140     the
141     child's children to the function.
142
143     called This is the number of times the function called
144     this child '/' the total number of times the child
145     was called. Recursive calls by the child are not
146     listed in the number after the '/'.
147
148     name This is the name of the child. The child's index
149     number is printed after it. If the child is a
150     member of a cycle, the cycle number is printed
151     between the name and the index number.
152
153 If there are any cycles (circles) in the call graph, there is an
154 entry for the cycle-as-a-whole. This entry shows who called the
155 cycle (as parents) and the members of the cycle (as children.)
156 The '+' recursive calls entry shows the number of function calls
157 that
158 were internal to the cycle, and the calls entry for each member
159 shows,
160 for that member, how many times it was called from other members
161 of
162 the cycle.
163
164 Copyright (C) 2012-2020 Free Software Foundation, Inc.
165
166 Copying and distribution of this file, with or without modification
167 ,
168 are permitted in any medium without royalty provided the copyright
169 notice and this notice are preserved.
170
171 Index by function name
172
173 [3] allocate                [2] finite_difference    [5]
174     potential
175 [4] allocate_complex        [1] main                [6]
176     wave_packet

```

8.2 Análise de resultados

Como esperado, no profile do código otimizado, a função que mais consome processamento ainda é a função que aplica o método das diferenças finitas.

Também é possível observar que, tanto para o profile sem e com as flags, a otimização nas outras funções melhorou o desempenho do código. Embora o efeito dessa otimização tenha um impacto mais significativo sem as melhores flags, obtendo uma redução no tempo de execução de aproximadamente seis minutos, continua sendo mais vantajoso utilizar as melhores flags mesmo que o seu tempo tenha reduzido em apenas vinte segundos após as otimizações.

Portanto, a otimização do código por meio da aplicação das boas práticas de programação é indispensável para melhorar o desempenho do programa serial.

No entanto, essa técnica não é suficiente para reduzir o tempo de execução drasticamente como proposto, sendo necessário recorrer à outros artifícios (no caso, ao paralelismo) para melhorar o desempenho do programa.

9 Implementação do programa com CUDA

A paralelização do programa será realizada por intermédio da interface de programação de aplicação CUDA para GPGPU criada pela Nvidia. As placas gráficas devem ser compatíveis com a API, isto é, geralmente essa devem possuir um chipset desenvolvido pela Nvidia. Abaixo está uma tabela com informações úteis sobre as placas gráficas utilizadas no projeto.

	CUDA cores	Arquitetura	Memória	DPFLOPS
Tesla K40t	2880	Volta	12GB	1.43
Tesla V100	5120	Volta	16-32GB	7.8

Tabela 1: Tabela com as informações de cada GPU a ser utilizada na implementação CUDA desse programa. Os valores para Double Precision FLOPS (DPFLOPS) são dados em TeraFLOPS.

Como é possível observar, a segunda placa gráfica possui uma quantidade maior de DPFLOPS, o que pode indicar um melhor desempenho em relação à primeira.

9.1 Estratégia de paralelização

A ideia central da paralelização desse problema concentra-se na conversão de uma matriz NxN por um vetor de tamanho equivalente. Deste modo, a estratégia adotada utiliza de N blocos com N threads cada, onde associa-se os blocos com as linhas da matriz e as threads com as colunas. Todos os blocos pertencem ao mesmo grid para que seja possível realizar a comunicação entre os dados de dependência temporal e a dimensão do bloco servirá para orientar a posição do bloco na grade.

```
1 matrix[i][j] -> matrix[i * size + j]
2 matrix[i*size + j] -> matrix[blockIdx.x * blockDim.x + threadIdx.x]
```

9.2 Profile

O profile para o código paralelizado com CUDA leva em consideração tanto a análise a nível de API quanto para a nível de GPU, uma vez que o comando utilizado foi:

```
1 nvprof -s -o schrodinger.nvprof ./schrodinger.x
```

Logo, segue as informações do profile.

```
[bianca.gomes@sduumont14 cuda]$ cat profile.txt
==28954== NVPROF is profiling process 28954, command: ./schrodinger.x
==28954== Error: File already exists: /scratch/uff21hpc/bianca.gomes/cuda/schrodinger.nvprof
==28954== Use "-f" to overwrite existing file
==28954== Profiling application: ./schrodinger.x
==28954== Profiling result:
   Type  Time(%)   Time    Calls   Avg       Min       Max   Name
GPU activities:  97.86%  36.1716s  19250  1.8790ms  1.6329ms  2.3578ms  [CUDA memcpy HtoD]
                2.14%  792.57ms  2750   288.21us  285.51us  292.55us  finite_difference_kernel
API calls:      99.27%  42.5558s  38500  1.1053ms  2.0310us  2.8927ms  cudaMemcpy
                0.39%  165.20ms   7    23.600ms  133.63us  164.36ms  cudaMalloc
                0.31%  133.46ms  2750   48.532us  42.264us  301.91us  cudaLaunchKernel
                0.02%  7.9675ms  19250   413ns    255ns    617.62us  cudaGetErrorString
                0.01%  4.3994ms   7    628.49us  326.04us  688.15us  cudaFree
                0.00%  912.90us   2    456.45us  446.80us  466.11us  cuDeviceTotalMem
                0.00%  542.22us   2    271.11us  250.15us  292.07us  cudaGetDeviceProperties
                0.00%  509.24us  194    2.6240us   147ns    129.18us  cuDeviceGetAttribute
                0.00%  49.971us   2    24.985us  23.498us  26.473us  cuDeviceGetName
                0.00%  14.285us   2    7.1420us  2.6400us  11.645us  cuDeviceGetPCIBusId
                0.00%  6.3310us   1    6.3310us  6.3310us  6.3310us  cudaGetDeviceCount
                0.00%  3.8690us   4     967ns    200ns    2.8780us  cuDeviceGet
                0.00%  1.9860us   3     662ns    160ns    1.1910us  cuDeviceGetCount
                0.00%    625ns    2     312ns    212ns    413ns    cuDeviceGetUuid
```

Figura 10: Profile para o código em CUDA. É fácil observar que a maior parte das atividades estão nas funções de cópia de *host* para *device* e vice-versa nomeadas cudaMemcpy.

10 Benchmark (paralelo)

Preliminarmente, é essencial a informação de que para as análises a seguir o tempo final e a matriz do programa foram reduzidos. Essa redução foi necessária por causa do tempo das filas de submissão dos clusters utilizados.

O tempo final será sempre

```
1 tempo final : 400
```

e o tamanho da matriz irá variar entre

```
1 cudacores="2880 5760 8640"
```

na finalidade de se ajustar aos cudacores em relação a GPU Tesla K40t e sem ultrapassar o tempo das filas de submissão.

10.1 Tesla K40t (LNCC)

Decorrente da sobrecarga de comunicação entre *host* e *device* e também da quantidade de DPFLOPS dessa placa gráfica, o tempo do programa paralelizado ultrapassa o tempo serial para todos os casos.

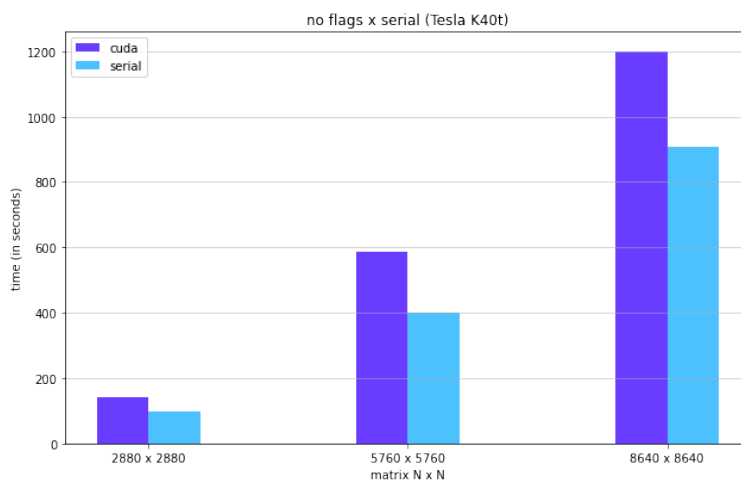


Figura 11: Benchmark para a comparação entre os tempos do programa serial e os tempos para o paralelizado com CUDA sem utilizar flags.

A eficiência é pior para o segundo caso onde há uma matriz 5760 x 5760, pode-se considerar que a relação entre o tamanho da malha e a quantidade de cudacores da GPU está prejudicando com processos de *pipeline*. De qualquer forma, para todos os casos, o programa é simplesmente ineficiente.

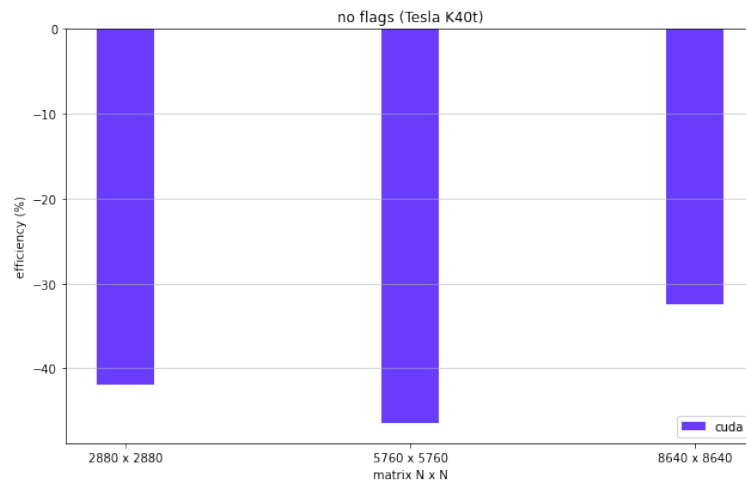


Figura 12: Benchmark com a eficiência do programa na GPU Tesla K40t em comparação com o tempo do programa serial. Ambos estão sem nenhuma flag.

O benchmark do speed up é completamente coerente para um programa ineficiente. Não há ganho em tempo, pelo contrário, é preferível executar o serial do que o paralelo. Enquanto que o parâmetro para o gráfico plano é a reta $x = y$, o parâmetro para o histograma é o valor um, portanto é razoável que todos os valores marcados sejam menores que esse número.

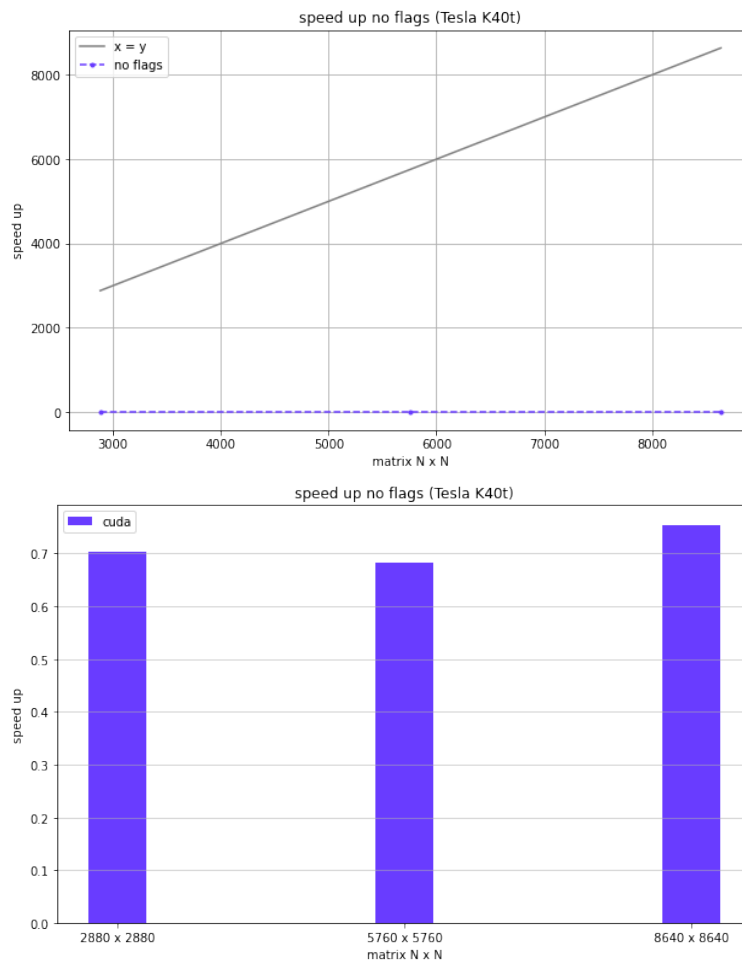


Figura 13: Benchmark para o speed up do programa paralelizado. Sem nenhuma flag.

10.1.1 Análise de flags

As flags utilizadas para aumentar a eficiência da implementação em CUDA estão dispostas abaixo. Todas estão acompanhadas de uma flag de otimização $-OX$ para X de 1 até 3.

```

1 declare -a flagscuda=(
2     "1" " "
3     "2" "-Xptxas"
4     "3" "-use_fast_math"
5 )

```

A segunda flag refere-se a aplicação das otimizações no *device*.

No primeiro caso, que é para a primeira configuração de matriz, as flags do *device* são as piores para todas as otimizações, infere-se que é por conta das operações que ocorrem no *kernel*. Não há grandes diferenças entre os tempos do primeiro e terceiro conjunto de flags, portanto considera-se a melhor flags sendo

`-O3 -use_fast_math`

nessa situação.

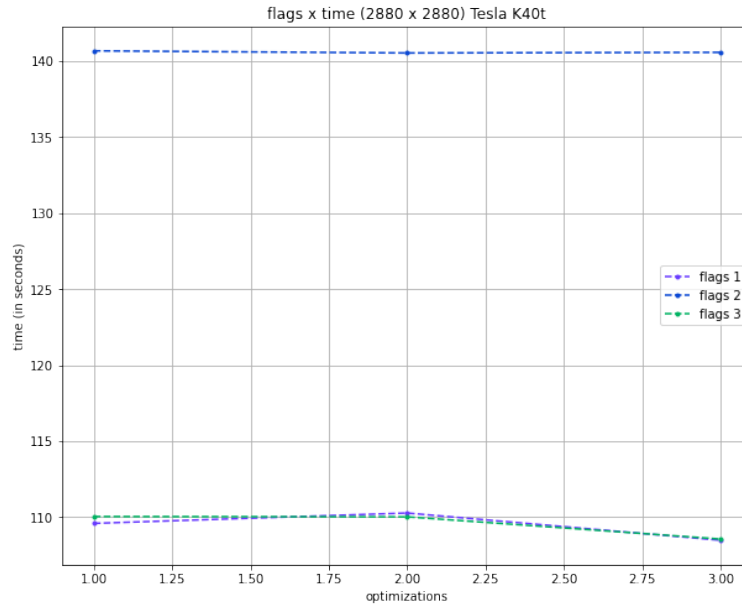


Figura 14: Benchmark do tempo para o conjunto de flags na matriz 2880 x 2880.

As melhores flags para a matriz 5760 x 5760 são as mesmas da análise anterior. No entanto, nota-se aqui um comportamento peculiar na otimização -O2, uma espécie de inversão de desempenho entre as flags de otimização do *host* e *device*. A melhor suposição é de que tal comportamento é decorrente do *pipeline*.

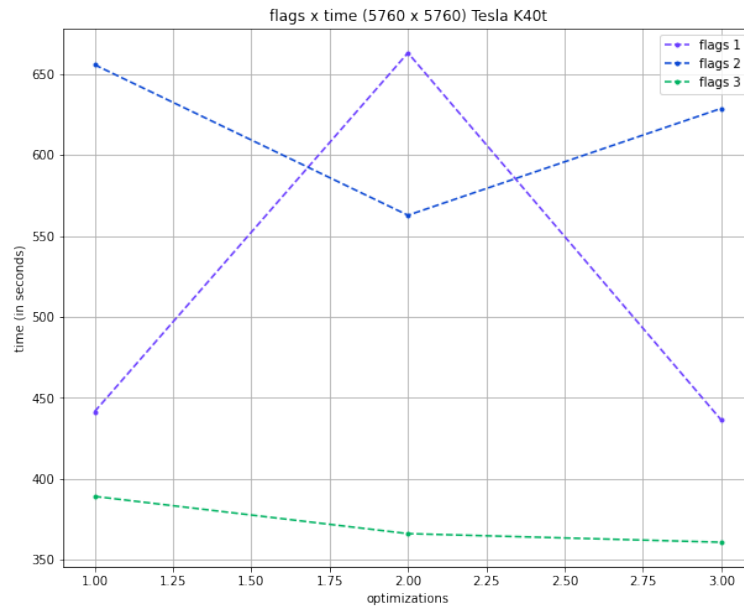


Figura 15: Benchmark do tempo para o conjunto de flags na matriz 5760 x 5760.

Novamente, as melhores flags são as mesmas das duas análises anteriores. O comportamento de inversão entre *host* e *device* aparece novamente.

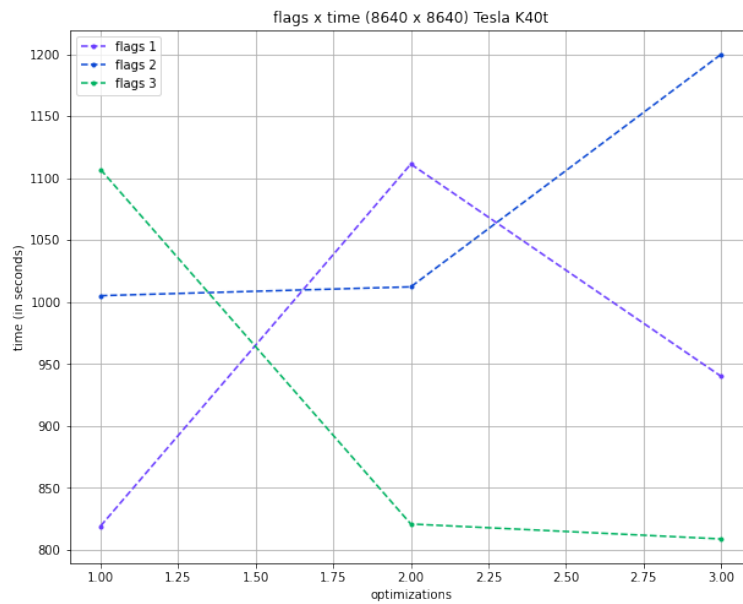


Figura 16: Benchmark do tempo para o conjunto de flags na matriz 8640 x 8640.

No geral, não há muito o que comentar sobre as eficiências separadamente pois essas se comportam de acordo com a análise temporal. O relevante aqui é enfatizar que, para a placa Tesla K40t, o programa com a matriz 5760 x 5760 unido das melhores flags atingiu a maior eficiência que está em torno de 40%.

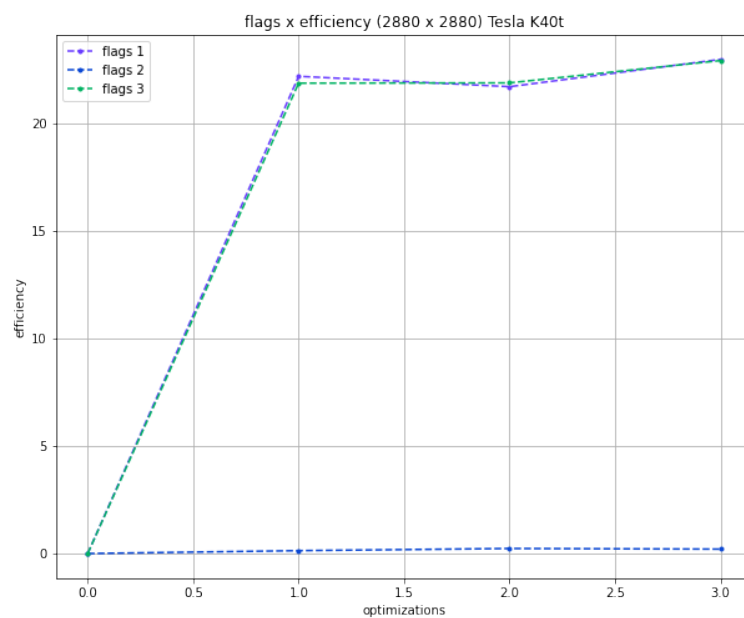


Figura 17: Benchmark da eficiência para o conjunto de flags na matriz 2880 x 2880.

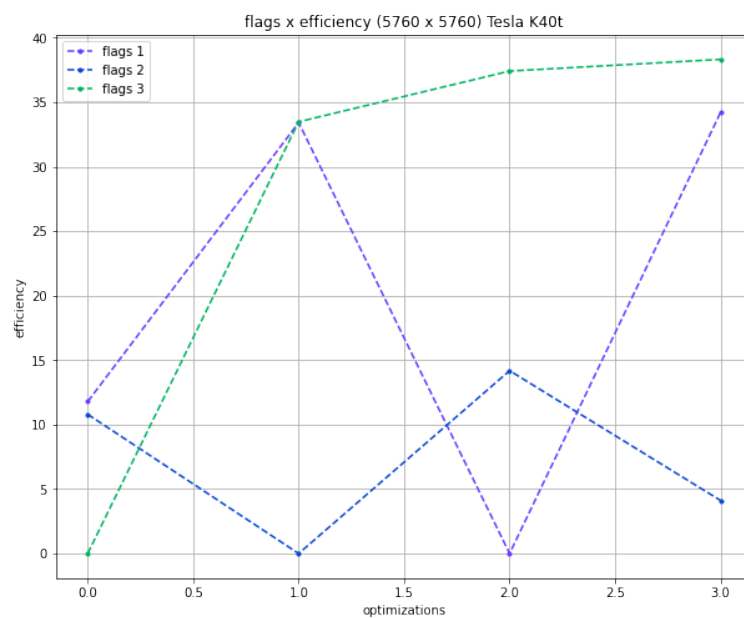


Figura 18: Benchmark da eficiência para o conjunto de flags na matriz 5760 x 5760.

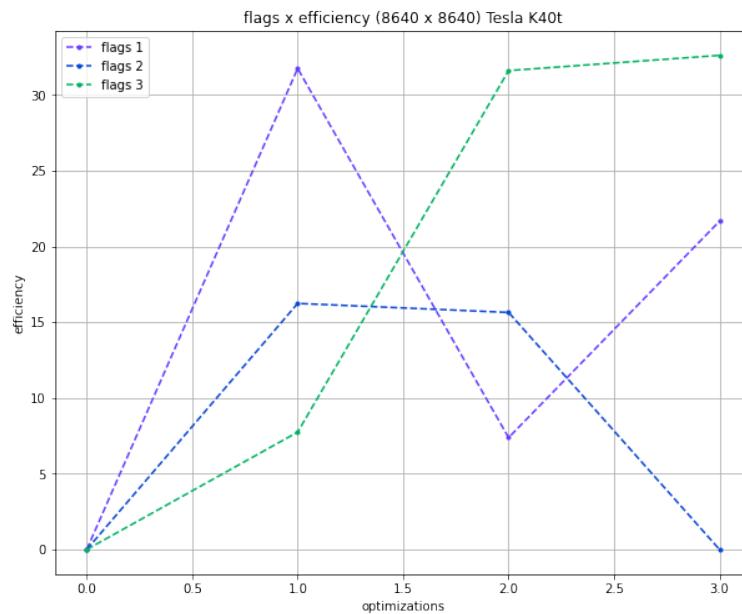


Figura 19: Benchmark da eficiência para o conjunto de flags na matrix 8640 x 8640.

10.2 Tesla V100 (LNCC)

Nessa placa gráfica, nota-se uma melhoria do tempo de execução do programa paralelo em relação ao serial. Esse melhor desempenho deve-se à quantidade de DPFLOPS presentes nessa placa.

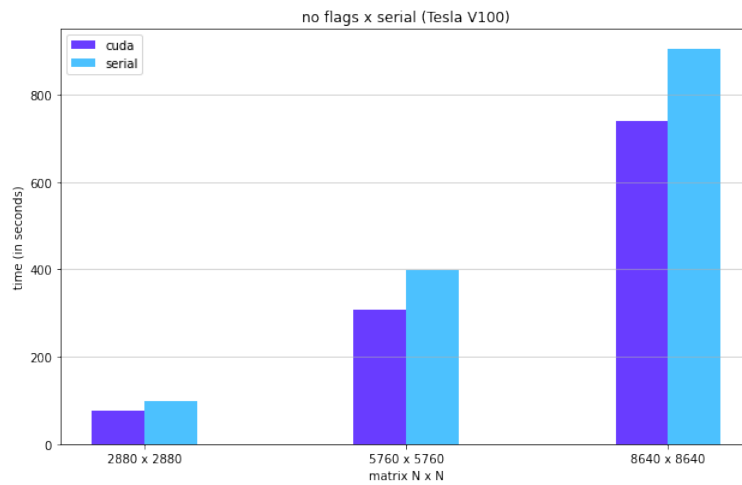


Figura 20: Benchmark dos tempos do programa paralelizado em comparação com os tempos do programa serial.

Ainda assim, a eficiência é considerada pouca, uma vez que não chega nem em 50%. As melhores eficiências são correspondentes às duas primeiras configurações de matriz, considera-se que isso acontece por causa da quantidade de CUDA cores na placa ser maior do que a primeira matriz e apenas um pouco menor do que a segunda, de tal modo que o *pipeline* não é tão agressivo para o desempenho do programa. No entanto, é nítido que o mesmo não aconteceu para a terceira configuração de matriz e para esse o programa é nitidamente menos eficiente.

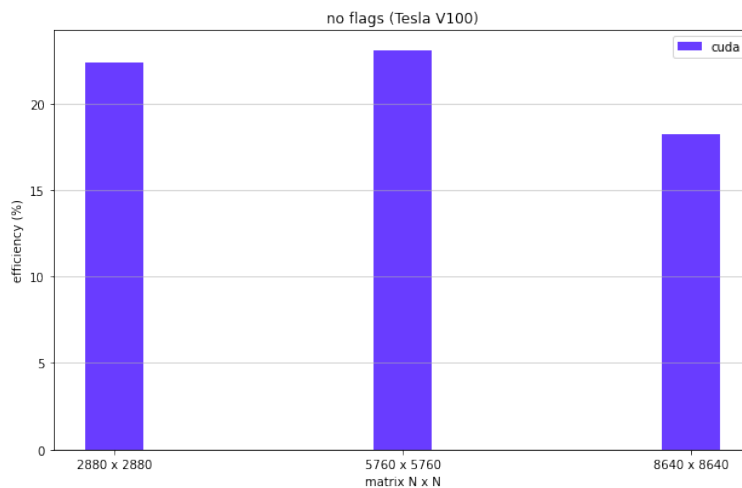


Figura 21: Benchmark da eficiência do programa paralelizado em CUDA.

Há ganho em eficiência no programa, mas é pequeno. Nota-se a distância entre o speed up e a reta de referência, como também nota-se que o aumento é decimal com relação ao número de referência do histograma.

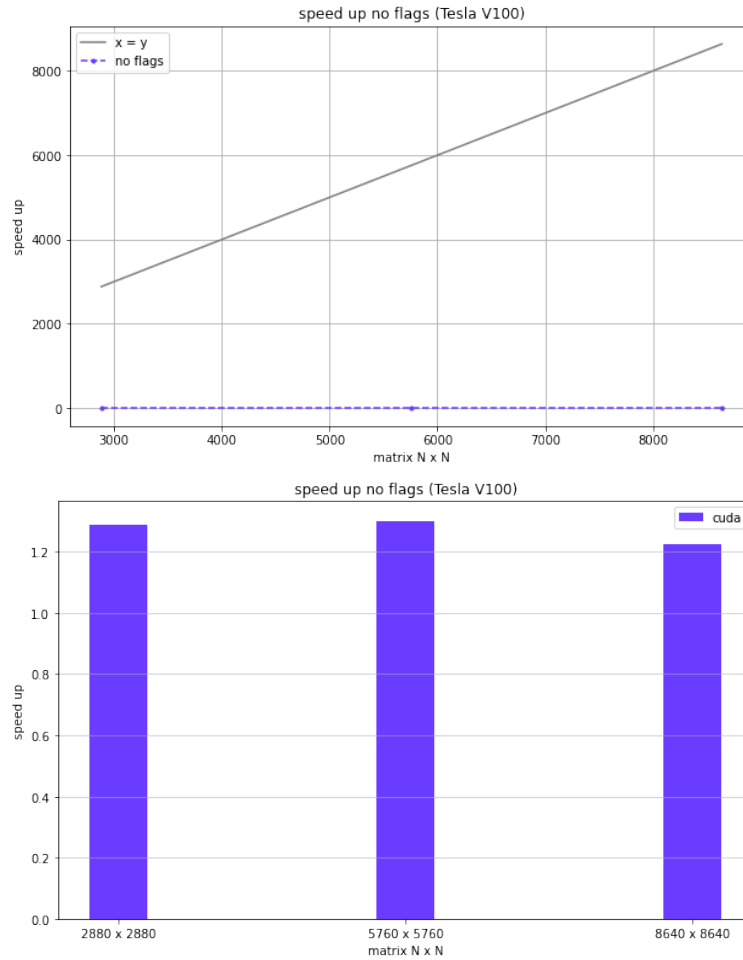


Figura 22: Benchmark do speed up do programa paralelizado executado com a placa Tesla V100.

10.2.1 Análise de flags

As mesmas flags da análise para a placa de vídeo anterior estão sendo utilizadas.

Semelhante à situação na outra placa gráfica, a flag do *device* é a pior para todas as otimizações. No entanto, aqui tem-se que a melhor flag é simplesmente a otimização -O3.

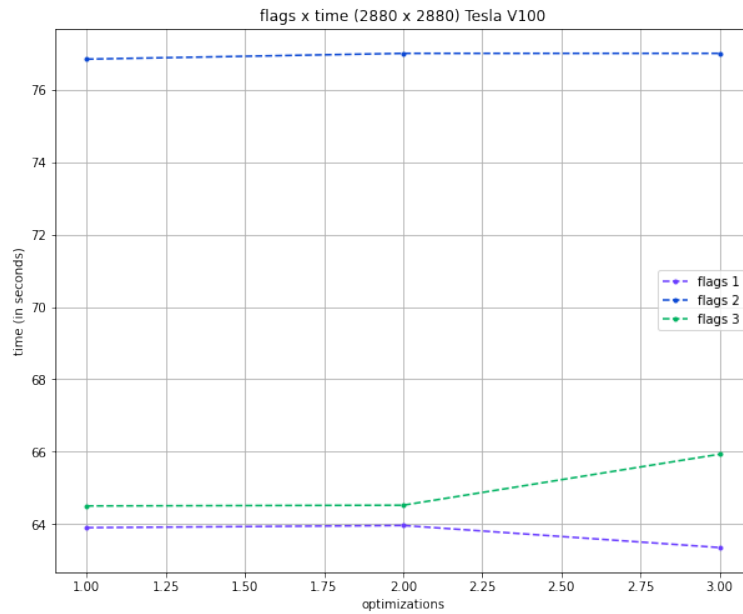


Figura 23: Benchmark do tempo para o conjunto de flags na matriz 2880 x 2880.

As melhores flags ficam entre as que carregam a otimização -O3 do primeiro e terceiro conjunto de flags. Consideraremos que a melhor é -O3 por causa da situação na matriz 2880 x 2880. Não há nenhum comportamento atípico nas flags do *device*, o que pode ser atribuído como consequência da quantidade de cudacores e DPFLOPS da placa gráfica.

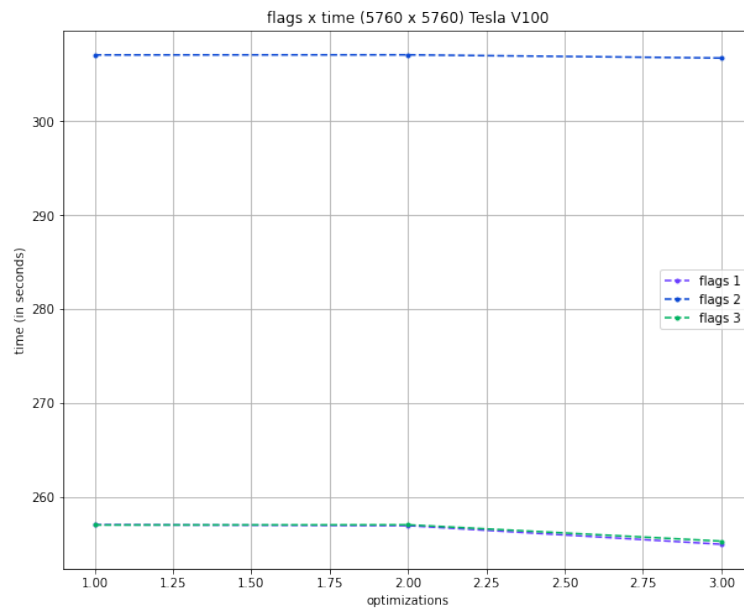


Figura 24: Benchmark do tempo para o conjunto de flags na matriz 5760 x 5760.

No entanto, para a matriz de tamanho 8640 x 8640, a melhor flag é

`-O3 -use_fast_math`

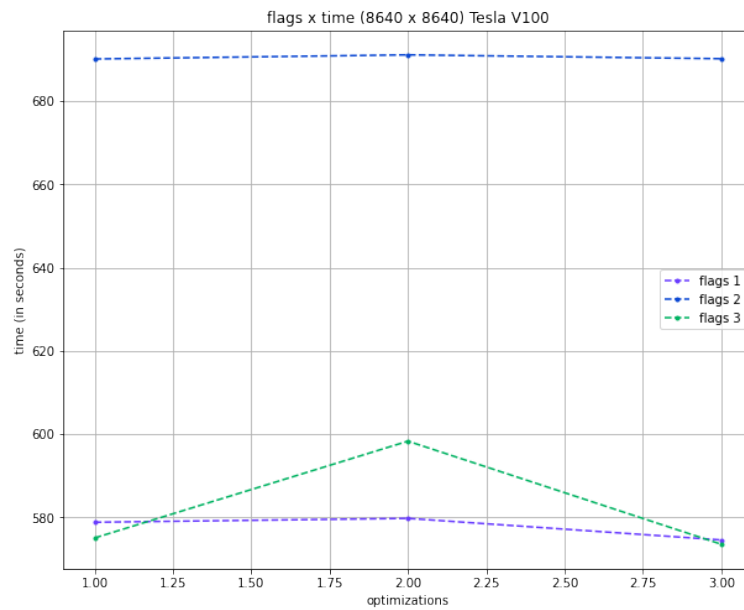


Figura 25: Benchmark do tempo para o conjunto de flags na matriz 8640 x 8640.

Novamente, como as eficiências seguem o comportamento do desempenho temporal, o importante a ser destacado é que a configuração de matriz mais eficiente é a 8640 x 8640.

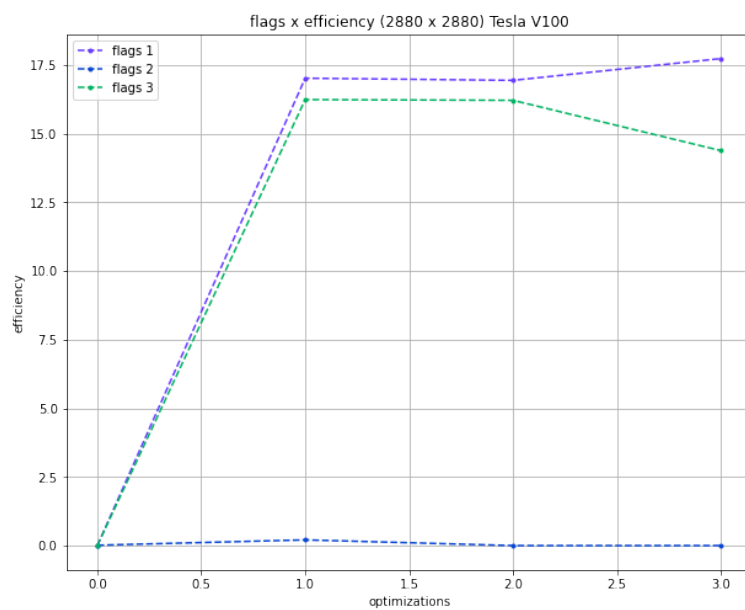


Figura 26: Benchmark da eficiência para o conjunto de flags na matriz 2880 x 2880.

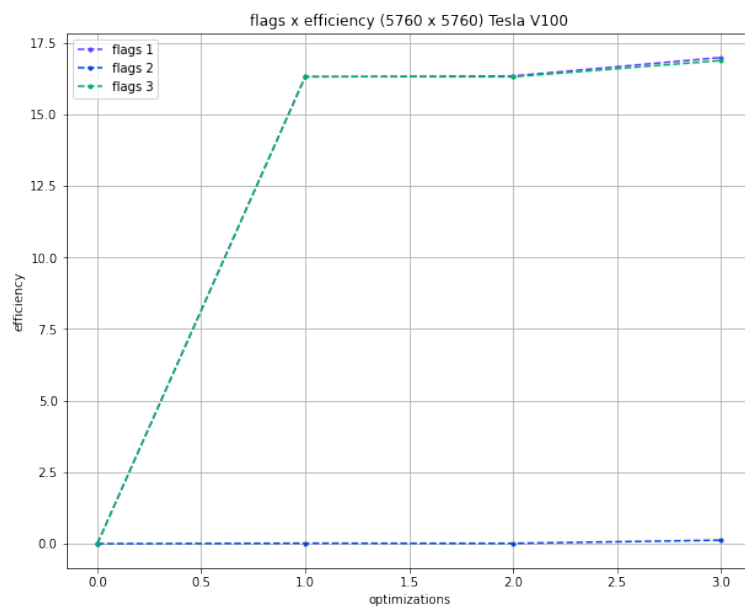


Figura 27: Benchmark da eficiência para o conjunto de flags na matriz 5760 x 5760.

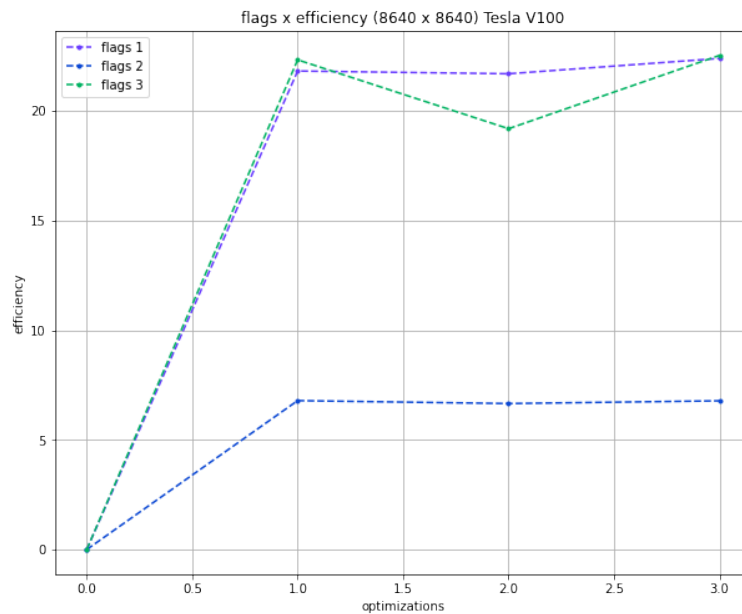


Figura 28: Benchmark da eficiência para o conjunto de flags na matriz 8640 x 8640.

10.3 Benchmark com as melhores flags (Tesla V100)

Dessa vez, com as melhores flags, os tempos obtidos para o código em CUDA são melhores do que os tempos do código em serial com uma diferença satisfatória no tempo de execução.

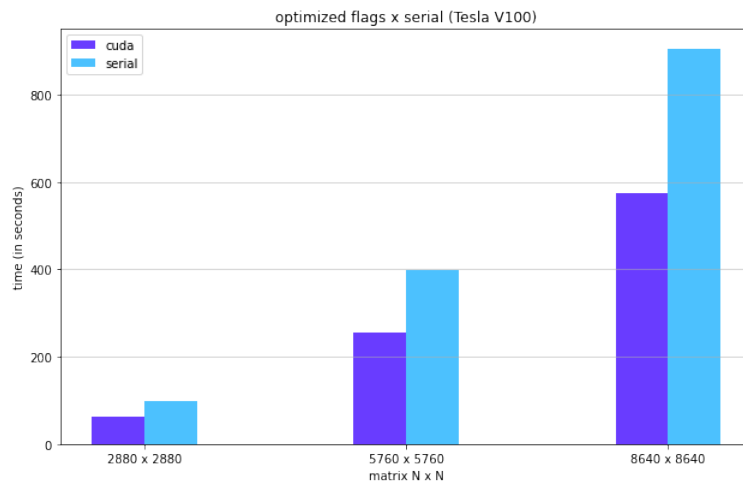


Figura 29: Benchmark da comparação entre os tempos otimizados com flags do código paralelo e os tempos do código serial.

A eficiência é nivelada apresentando valores semelhantes para todos os conjuntos de matrizes analisados.

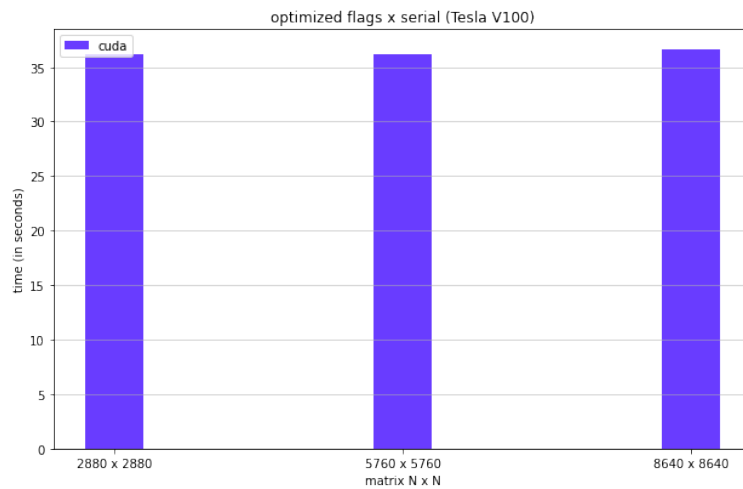


Figura 30: Benchmark da comparação entre as eficiências otimizadas com flags do código paralelo e os tempos do código serial.

11 Comparação com a implementação em OpenMP

O benchmark abaixo demonstra a evolução da eficiência do programa que anteriormente foi paralelizado com OpenMP. A execução desse programa foi realizado em um computador caseiro que tem as mesmas configurações que o computador do Laboratório 107C da UFF.

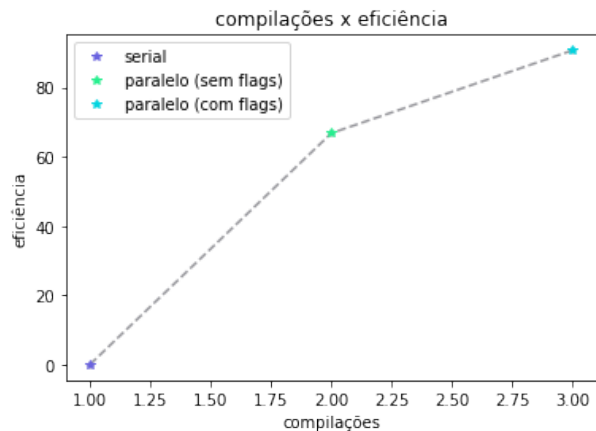


Figura 31: Gráfico para o desenvolvimento da eficiência por compilação do programa paralelizado com OpenMP.

O benchmark para o programa paralelizado em CUDA leva em consideração a matriz da melhor eficiência obtida com o uso de flags.

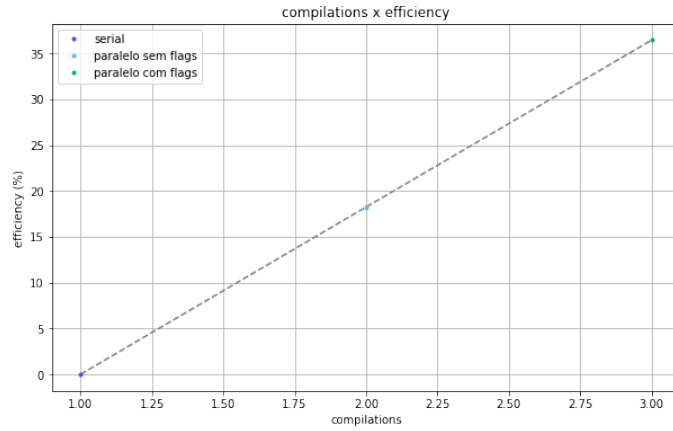


Figura 32: Gráfico para o desenvolvimento da eficiência por compilação do programa paralelizado com CUDA. A matriz utilizada foi a de tamanho 8640 x 8640 pois essa obteve a melhor eficiência com flags. Os dados referem-se a execução do programa na placa Tesla V100.

A conclusão de que o programa em OpenMP é superior em questão de desempenho temporal é imediato, pois enquanto esse obteve uma eficiência em torno de 80 – 90% em um computador caseiro, o programa em CUDA obteve em torno de 35% de eficiência em uma placa gráfica voltada para computação de alto desempenho.

11.1 Considerações sobre o método das diferenças finitas

É evidente pelos resultados obtidos nos benchmarks que essa estratégia de paralelização adotada é ineficiente ou pouco eficiente.

O método das diferenças finitas como foi aplicado nessa solução para o problema da equação de Schrödinger, obrigatoriamente exige a separação entre as partes reais e imaginárias da equação. Consequentemente, isso sobrecarrega a quantidade de informação que deve ser comunicada entre o *host* e o *device*, que é possível confirmar a ocorrência ao observar o profile do código.

Uma alternativa em potencial para resolver esse problema de eficiência é substituir o método número vigente por algum método que utilize da *Fast Fourier Transform*. Obviamente, não descarta-se a possibilidade de outra estratégia de paralelização com o método atual que seja mais eficiente. No entanto, a eficácia da FFT com aplicações em CUDA é reconhecida, existindo até mesmo uma *library* da Nvidia para isso [3].

12 Validação de Resultados

Dos postulados da mecânica quântica tem-se a declaração essencial de que a equação de Schrödinger é responsável pela evolução temporal do sistema. Com isso, é equivalente dizer, tendo em mente a interpretação probabilística, que a normalização da função de onda se conserva com o tempo. Portanto, para validar os resultados da implementação computacional para a solução dessa equação, considera-se a densidade de probabilidade descrita pela equação abaixo.

$$P(x, y, t)dx dy = |\psi(x, y, t)|^2 dx dy \quad (9)$$

Com os resultados de compilações do programa para uma matriz reduzida e tempos finais diferentes, obteve-se as seguintes imagens:

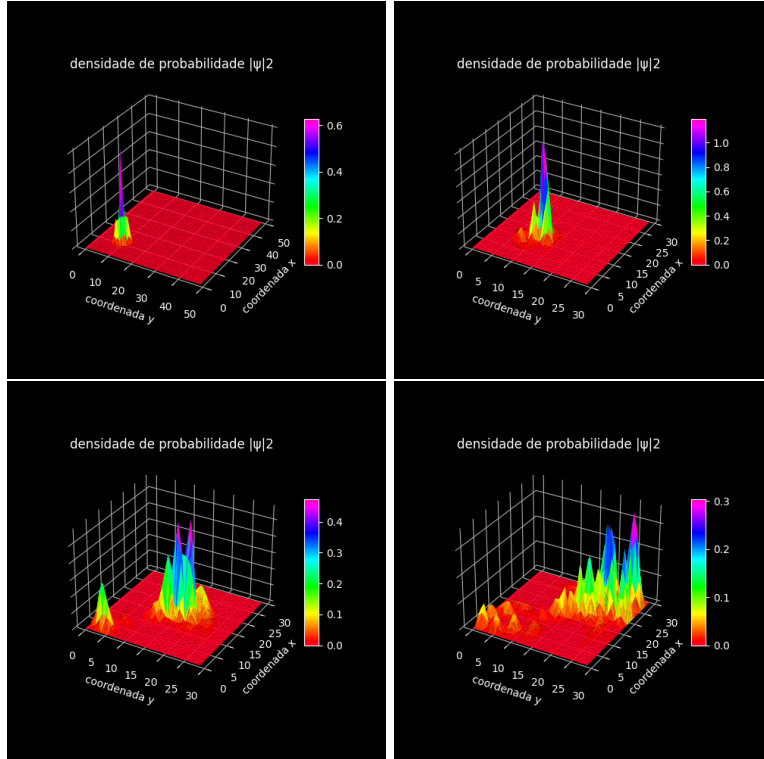


Figura 33: Representação gráfica obtida com a compilação do programa serial para tempos finais diferentes, respectivamente, $t_f=0$, $t_f=100000$, $t_f=500000$ e $t_f=1000000$. Nesse caso, a matriz está reduzida.

Com essas imagens nota-se que a função de onda transita pelo potencial conforme ocorre a evolução temporal do sistema. Outra importantíssima

característica é a preservação da distribuição gaussiana, tanto para a parte transmitida quanto para a parte refletida, conservando a interpretação probabilística.

Para confirmar a preservação dos dados após a paralelização do programa e também para que não seja necessário uma avaliação trabalhosa de dado por dado, ainda considerando a interpretação probabilística do problema tem-se como parâmetro a média e o desvio padrão dos dados.

$$\mu = \frac{\sum_i^n x_i}{n} \quad (10)$$

$$\sigma = \sqrt{\frac{\sum_i^n (x_i - \bar{x})^2}{n}} \quad (11)$$

Implementação	Média μ	Desvio Padrão σ
Serial	$8.00 \cdot 10^{-8}$	$1.78 \cdot 10^{-5}$
CUDA	$1.56 \cdot 10^{-7}$	$2.80 \cdot 10^{-5}$

Tabela 2: Tabela com os valores de média e desvio padrão para cada implementação do programa.

O critério de qualidade desses dados é definido pela discrepância dos dados do código em paralelo com os dados do código em serial.

	Discrepância da média	Discrepância do desvio padrão
CUDA	$7.65 \cdot 10^{-8}$	$1.02 \cdot 10^{-5}$

Tabela 3: Tabela com as discrepâncias entre os dados do código paralelizado e sua versão serial.

Como as discrepâncias entre os dois conjuntos de dados são extremamente pequenas, pode-se considerar que a implementação da paralelização em CUDA não alterou a qualidade dos dados e que a solução do problema foi efetivada com sucesso.

13 Conclusões

Com o auxílio das melhores flags obtidas com os benchmarks para o programa serial realizou-se o profile, onde ficou claro que a função responsável pelo método das diferenças finitas consumia mais processamento que todas as outras e que a paralelização deveria atuar diretamente nessa função e no loop onde a mesma é chamada.

O método das diferenças finitas não aparenta ser o melhor método para receber uma implementação CUDA do problema abordado, dado que esse exige a separação entre as partes real e imaginária da equação de onda quântica e portanto, sobrecarregando a comunicação entre *host* e *device*. Ainda assim, com a placa Tesla V100 obteve-se eficiência por mais que seja pouca. No entanto, com a placa Tesla K40t, o programa foi simplesmente ineficiente inferindo que será sempre necessário uma placa com foco em computação de alto desempenho para que o programa seja minimamente eficiente.

Finalmente, diante da validação dos resultados, percebendo que as discrepâncias entre as médias e desvios padrões do programa serial e do programa paralelizado são extremamente pequenas, pode-se concluir que a paralelização foi concluída corretamente, embora os ganhos em desempenho temporal não sejam satisfatórios.

14 Apêndice

14.1 Informações sobre as CPUs

14.1.1 Laboratório 107C (UFF)

```
1 Architecture:          x86_64
2 CPU op-mode(s):       32-bit, 64-bit
3 Byte Order:           Little Endian
4 CPU(s):                8
5 On-line CPU(s) list:  0-7
6 Thread(s) per core:   2
7 Core(s) per socket:   4
8 Socket(s):            1
9 NUMA node(s):         1
10 Vendor ID:            GenuineIntel
11 CPU family:           6
12 Model:                42
13 Model name:           Intel(R) Core(TM) i7-2600 CPU @ 3.40GHz
14 Stepping:             7
15 CPU MHz:              1599.975
16 CPU max MHz:          3800,0000
17 CPU min MHz:          1600,0000
18 BogoMIPS:             6784.74
19 Virtualization:       VT-x
20 L1d cache:            32K
21 L1i cache:            32K
22 L2 cache:             256K
23 L3 cache:             8192K
24 NUMA node0 CPU(s):    0-7
25 Flags:                fpu vme de pse tsc msr pae mce cx8 apic
26 sep mtrr pge mca cmov pat pse36 clflush dts acpi mmx fxsr sse
27 sse2 ssht tm pbe syscall nx rdtscp lm constant_tsc arch_perfmon
28 pebs bts rep_good nopl xtopology nonstop_tsc aperfmperf
29 eagerfpu pni pclmulqdq dtes64 monitor ds_cpl vmx smx est tm2
30 ssse3 cx16 xtpr pdcm pcid sse4_1 sse4_2 x2apic popcnt
31 tsc_deadline_timer aes xsave avx lahf_lm epb ssbd ibrs ibpb
32 stibp tpr_shadow vnmi flexpriority ept vpid xsaveopt dtherm
33 ida arat pln pts md_clear spec_ctrl intel_stibp flush_l1d
```

14.1.2 B107 (LNCC)

```
1 Architecture:          x86_64
2 CPU op-mode(s):       32-bit, 64-bit
3 Byte Order:           Little Endian
4 CPU(s):                24
5 On-line CPU(s) list:  0-23
6 Thread(s) per core:   1
7 Core(s) per socket:   12
8 Socket(s):            2
9 NUMA node(s):         2
10 Vendor ID:            GenuineIntel
11 CPU family:           6
12 Model:                62
13 Model name:           Intel(R) Xeon(R) CPU E5-2695 v2 @ 2.40GHz
```

```

14 Stepping: 4
15 CPU MHz: 2899.072
16 CPU max MHz: 3200,0000
17 CPU min MHz: 1200,0000
18 BogoMIPS: 4800.35
19 Virtualization: VT-x
20 L1d cache: 32K
21 L1i cache: 32K
22 L2 cache: 256K
23 L3 cache: 30720K
24 NUMA node0 CPU(s): 0-11
25 NUMA node1 CPU(s): 12-23
26 Flags: fpu vme de pse tsc msr pae mce cx8 apic
27 sep mtrr pge mca cmov pat pse36 clflush dts acpi mmx fxsr sse
28 sse2 ss ht tm pbe syscall nx pdpe1gb rdtscp lm constant_tsc
29 arch_perfmon pebs bts rep_good nopl xtopology nonstop_tsc
30 aperfmperf eagerfpu pni pclmulqdq dtes64 monitor ds_cpl vmx
31 smx est tm2 ssse3 cx16 xtpr pdcm pcid dca sse4_1 sse4_2
32 x2apic popcnt tsc_deadline_timer aes xsave avx f16c rdrand
33 lahf_lm epb ssbd ibrs ibpb stibp tpr_shadow vnmi flexpriority
34 ept vpid fsgsbase smep erms xsaveopt dtherm ida arat pln pts
35 md_clear spec_ctrl intel_stibp flush_l1d

```

14.1.3 SequanaX (LNCC)

```

1 Architecture: x86_64
2 CPU op-mode(s): 32-bit, 64-bit
3 Byte Order: Little Endian
4 CPU(s): 96
5 On-line CPU(s) list: 0-47
6 Off-line CPU(s) list: 48-95
7 Thread(s) per core: 1
8 Core(s) per socket: 24
9 Socket(s): 2
10 NUMA node(s): 2
11 Vendor ID: GenuineIntel
12 CPU family: 6
13 Model: 85
14 Model name: Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz
15 Stepping: 7
16 CPU MHz: 2101.000
17 CPU max MHz: 2101,0000
18 CPU min MHz: 1000,0000
19 BogoMIPS: 4200.00
20 Virtualization: VT-x
21 L1d cache: 32K
22 L1i cache: 32K
23 L2 cache: 1024K
24 L3 cache: 36608K
25 NUMA node0 CPU(s): 0-23
26 NUMA node1 CPU(s): 24-47
27 Flags: fpu vme de pse tsc msr pae mce cx8 apic
28 sep mtrr pge mca cmov pat pse36 clflush dts acpi mmx fxsr sse
29 sse2 ss ht tm pbe syscall nx pdpe1gb rdtscp lm constant_tsc
30 art arch_perfmon pebs bts rep_good nopl xtopology nonstop_tsc
31 aperfmperf eagerfpu pni pclmulqdq dtes64 monitor ds_cpl vmx

```

```

32 smx est tm2 ssse3 sdbg fma cx16 xtpr pdcm pcid dca sse4_1
33 sse4_2 x2apic movbe popcnt tsc_deadline_timer aes xsave avx
34 f16c rdrand lahf_lm abm 3dnowprefetch epb cat_l3 cdp_l3
35 invpcid_single intel_ppin intel_pt ssbd mba ibrs ibpb stibp
36 ibrs_enhanced tpr_shadow vnmi flexpriority ept vpid fsgsbase
37 tsc_adjust bmi1 hle avx2 smep bmi2 erms invpcid rtm cqm mpx
38 rdt_a avx512f avx512dq rdseed adx clflushopt clwb avx512cd
39 avx512bw avx512vlxsavopt xsavec xgetbv1 cqm_llc cqm_occup_llc
40 cqm_mbm_total cqm_mbm_local dtherm ida arat pln pts hwp
41 hwp_act_window hwp_epp hwp_pkg_req pku ospke avx512_vnni
42 md_clear spec_ctrl intel_stibp flush_l1d arch_capabilities

```

14.2 Informações sobre as GPUs

14.2.1 Tesla K40t

```

1 device 0 information
2     name: Tesla K40t
3     compute capability: 3.5
4     clock rate: 745000
5     device copy overlap: enabled
6
7     kernel execution timeout:
8     disabled
9
10 device memory information
11     total global memory: 11997020160
12     total constant memory: 65536
13     max memory pitch: 2147483647
14     texture alignment: 512
15
16 device multiprocessor information
17     multiprocessor count: 15
18     shared memory per multiprocessor: 49152
19     register per multiprocessor: 65536
20     threads in warp: 32
21     max thread dimensions: (1024, 1024, 64)
22     max grid dimensions: (2147483647, 65535, 65535)
23
24
25 device 1 information
26     name: Tesla K40t
27     compute capability: 3.5
28     clock rate: 745000
29     device copy overlap: enabled
30
31     kernel execution timeout:
32     disabled
33
34 device memory information
35     total global memory: 11997020160
36     total constant memory: 65536
37     max memory pitch: 2147483647
38     texture alignment: 512
39
40 device multiprocessor information

```

```

41     multiprocessor count: 15
42     shared memory per multiprocessor: 49152
43     register per multiprocessor: 65536
44     threads in warp: 32
45     max thread dimensions: (1024, 1024, 64)
46     max grid dimensions: (2147483647, 65535, 65535)

```

14.2.2 Tesla V100

```

1  device 0 information
2      name: Tesla V100-SXM2-32GB
3      compute capability: 7.0
4      clock rate: 1530000
5      device copy overlap: enabled
6
7      kernel execution timeout:
8      disabled
9
10 device memory information
11     total global memory: 34089730048
12     total constant memory: 65536
13     max memory pitch: 2147483647
14     texture alignment: 512
15
16 device multiprocessor information
17     multiprocessor count: 80
18     shared memory per multiprocessor: 49152
19     register per multiprocessor: 65536
20     threads in warp: 32
21     max thread dimensions: (1024, 1024, 64)
22     max grid dimensions: (2147483647, 65535, 65535)
23
24
25 device 1 information
26     name: Tesla V100-SXM2-32GB
27     compute capability: 7.0
28     clock rate: 1530000
29     device copy overlap: enabled
30
31     kernel execution timeout:
32     disabled
33
34 device memory information
35     total global memory: 34089730048
36     total constant memory: 65536
37     max memory pitch: 2147483647
38     texture alignment: 512
39
40 device multiprocessor information
41     multiprocessor count: 80
42     shared memory per multiprocessor: 49152
43     register per multiprocessor: 65536
44     threads in warp: 32
45     max thread dimensions: (1024, 1024, 64)
46     max grid dimensions: (2147483647, 65535, 65535)
47

```

```

48
49 device 2 information
50     name: Tesla V100-SXM2-32GB
51     compute capability: 7.0
52     clock rate: 1530000
53     device copy overlap: enabled
54
55     kernel execution timeout:
56     disabled
57
58 device memory information
59     total global memory: 34089730048
60     total constant memory: 65536
61     max memory pitch: 2147483647
62     texture alignment: 512
63
64 device multiprocessor information
65     multiprocessor count: 80
66     shared memory per multiprocessor: 49152
67     register per multiprocessor: 65536
68     threads in warp: 32
69     max thread dimensions: (1024, 1024, 64)
70     max grid dimensions: (2147483647, 65535, 65535)
71
72
73 device 3 information
74     name: Tesla V100-SXM2-32GB
75     compute capability: 7.0
76     clock rate: 1530000
77     device copy overlap: enabled
78
79     kernel execution timeout:
80     disabled
81
82 device memory information
83     total global memory: 34089730048
84     total constant memory: 65536
85     max memory pitch: 2147483647
86     texture alignment: 512
87
88 device multiprocessor information
89     multiprocessor count: 80
90     shared memory per multiprocessor: 49152
91     register per multiprocessor: 65536
92     threads in warp: 32
93     max thread dimensions: (1024, 1024, 64)
94     max grid dimensions: (2147483647, 65535, 65535)

```

Referências

- [1] Jos Thijssen. *Computational physics*. Cambridge university press, 2007. Citado na página 6.
- [2] H Rubin and Author Landau. *Computational Physics: Problem Solving with Computers*. J. Wiley & sons, 1997. Citado na página 6.
- [3] Fast fourier transforms for nvidia gpus. <https://developer.nvidia.com/cufft>. Citado na página 46.